

A Broadened Approach to Improved Estimation in Survey Sampling

Kyle Vincent and Christopher S. Henry *

Abstract

We outline a new framework for design-based inference in a survey sampling setting; we consider how estimation changes when a strategy (that is, a pairing of a sampling design and estimator) is first selected at random amongst a set of candidate strategies. Extending on the traditional approach in survey sampling, we redefine the observed data to now include a mention of the strategy used. The minimal sufficient statistic for the population parameter vector is derived for this setup; data reduction now involves considering hypothetical sampling designs that could have given rise to the sample along with their accompanying estimator(s). Implications on Rao-Blackwellization are investigated. Two studies are detailed to highlight how the new inference procedure can offer alternative views to, and circumvent some challenges encountered with, the traditional approach to inference.

Keywords: Design-based inference; Rao-Blackwell theorem; Minimal Sufficient statistic; Sampling strategy.

Acknowledgements

All views expressed in this manuscript are solely those of the authors and should not be attributed to the Bank of Canada.

*Currency Department, Bank of Canada, 234 Wellington Street, Ottawa, Ontario, CANADA, K1A 0G9
email: kvincent@bankofcanada.ca, chenry@bankofcanada.ca

1 Introduction

Traditional design-based inference in survey studies, also known as the fixed population approach, is based on an ordered set of observed sample unit labels together with their corresponding responses of interest. Well-known sources, like those of Cassel et al. (1977) and Thompson and Seber (1996), will denote the observed sample data as $d_0 = ((i, y_i) : i \in \underline{s})$, where \underline{s} is the set of units in the original order they were selected and y_i is the response of unit i . One may notice that there is no mention of the sampling design in the observed data, nor is there any mention of details regarding the estimation procedure. In this manuscript we show that if the strategy is chosen with some probability from amongst a set of candidate strategies then this results in a minimal sufficient statistic based on more than just the population-level data. The new setup to inference results in an approach that has potential to overcome some limitations that are present with the traditional approach, primarily because Rao-Blackwellization now entails averaging over estimators corresponding with multiple strategies.

We consider two cases where the new approach can lend itself well for estimation. First, Cassel et al. (1977) reviewed the admissibility properties of strategies and noted how no strategy is guaranteed to always perform better than any other strategy. Consequently, at times it may be unclear as to which strategy should serve best for a study. Further, some well-known strategies can give rise to highly-skewed sampling distributions for commonly used estimators, like when one is sampling from highly uneven populations with probability proportional to size depending on auxiliary information and the Hansen-Hurwitz estimator is used. Upon selecting such samples the analyst may feel tempted or inclined to make inference based on a different strategy. With our setup, we show how the analyst can exploit the corresponding Rao-Blackwell features to make this possible, which in turn can dampen the impact of designs on estimators sensitive to such possibly unrepresentative samples. Hence, with our method one can properly base inference on a conservative approach through taking a composite average of several strategies to avoid committing to one potentially weak strategy.

Second, we consider how the new setup can aid in a retrospective style of analysis. We consider situations when only a limited amount of information is available on the actual sampling design that was used, possibly because the procedure is beyond the full control of the analyst. We show how the analyst can utilize such information in the analysis with the aid of the new setup we present. The method is interesting in that it parallels the Bayesian paradigm as it allows the analyst to draw on prior knowledge of the design used to select the sample, as reflected upon in the discussion section of this manuscript.

Our goal is to offer a new and alternative view to design-based inference with ideas as to where the approach may be practical in an empirical setting. We consider empirical data

sets and simulation studies to justify our claims and encourage future work based on the new setup. The article is organized as follows. Section 2 introduces the setup with mathematical notation. Section 3 provides the minimal sufficiency result and a proof that the minimal sufficient statistic is not complete. Section 4 gives details and formulas oriented about the improved estimation procedure, which includes the variance expression and a corresponding expression for an estimate of the variance of the improved estimator. Section 5 provides several ideas with simulation results to demonstrate how the new setup is practical for some specific cases. Finally, Section 6 discusses the implications of the novel inference procedure outlined in this manuscript.

2 Framework and Mathematical Notation

Define $U = \{1, 2, \dots, N\}$ to be the population unit labels where N is the size of the population. Define $\underline{\theta} = (y_1, y_2, \dots, y_N)$ to be the population parameter vector; y_i is the response of unit i , $\underline{\theta} \in \Theta$ where $\Theta = \mathcal{R}^N$, and \mathcal{R} is the set of all real numbers.

Define \mathcal{D} to be the set of candidate sampling designs under consideration for the data collection aspect of the study. For ease of presentation we will assume that there are a finite number of sampling designs. Hence, the set of candidate designs can be ordered such that $\mathcal{D} = (\delta_1, \delta_2, \dots, \delta_T)$ where δ_k refers to sampling design k , $k = 1, 2, \dots, T$. For example, δ_1 may refer to the simple random sampling design and δ_2 may refer to a probability proportional to size sampling design.

Suppose $\phi = \phi(\underline{\theta})$ is a population quantity to be estimated. Define Φ_k to be the set of candidate estimators that accompany sampling design k . For ease of presentation we will assume that there are a finite number of estimators accompanying Φ_k . Hence, the corresponding set of candidate estimators can be ordered such that $\Phi_k = (\hat{\phi}_{k,1}, \hat{\phi}_{k,2}, \dots, \hat{\phi}_{k,V_k})$ where $\hat{\phi}_{k,j}$ refers to estimator j that accompanies sampling design k , $j = 1, 2, \dots, V_k$. For example, if δ_1 is the simple random sampling design then $\hat{\phi}_{1,1}$ may refer to the sample mean estimator and $\hat{\phi}_{1,2}$ may refer to the sample median estimator.

Define $\Lambda_{k,j}$ to be the strategy that corresponds with sampling design δ_k and estimator $\hat{\phi}_{k,j}$. Let $\underline{P}_\Lambda = (P(\Lambda_{k,j}) : k = 1, 2, \dots, T, j = 1, 2, \dots, V_k)$ be the probability vector with $P(\Lambda_{k,j})$ referring to the probability that strategy $\Lambda_{k,j}$ is selected for data collection and inference.

Define \underline{s} to be the vector of sample unit labels displayed in the order they were selected (with possible repeat selections). Define $P(\underline{s}|\underline{y}_s, \delta_k)$ to be the probability that \underline{s} is selected under sampling design k where \underline{y}_s is the corresponding vector of observed responses for all individuals

selected for the original sample \underline{s} .

The original data that is observed is $d_0 = (\Lambda_{k,j}, (i, y_i) : i \in \underline{s})$. Define r_d to be the reduction function that maps the original data to the reduced data $d_R = \{(i, y_i) : i \in s\}$ (that is $r_d(d_0) = d_R$) where s is the unordered set of unit labels. Define the reduced data ordered from smallest to largest label in the sample to be $d = (s, \underline{y}_s)$. As is common in the framework of design-based inference (see Thompson and Seber (1996) for a comprehensive review of the design-based setup to inference in survey sampling), we define a parameter vector $\underline{\theta}$ to be *consistent* with d if the i^{th} component of \underline{y}_s is equal to the s_i^{th} component of $\underline{\theta}$, say. We define Θ_d to be the set of all $\underline{\theta}$ that are consistent with d . Notice that, since consistency depends only on the values of the distinct units in the sample, $\Theta_d = \Theta_{d_R}$.

3 Minimal Sufficiency Result and Incompleteness

In the traditional setup the minimal sufficient statistic for $\underline{\theta} = (y_1, y_2, \dots, y_N)$ has been shown to be $d_R = \{(i, y_i) : i \in s\}$ (Pathak (1964); Godambe (1966); Basu (1969)). Cassel et al. (1977) later showed that the minimal sufficient is not complete. In this section we show that the minimal sufficient statistic in our inferential setup is still the aforementioned reduced data and that it is also not complete.

Theorem: Suppose that \underline{P}_Λ does not depend on $\underline{\theta}$ and that all candidate sampling designs are of the adaptive and/or conventional type (see Thompson and Seber (1996) for a discussion on such designs). That is, for all \underline{s} and $\Lambda_{k,j}$, $P(\underline{s}|\underline{y}, \Lambda_{k,j}) = P(\underline{s}|\underline{y}_s, \delta_k) (= P(\underline{s}|\delta_k))$ if the design is conventional). The random variable D_R is the minimal sufficient statistic for $\underline{\theta}$.

Proof: Choose any possible d_0 and consider Θ_{d_R} . Suppose that strategy $\Lambda_{k,j}$ is the strategy selected for the data collection and inference aspects of the study, for some arbitrary k and j . Now,

$$\begin{aligned} P_{\underline{\theta}}(D_0 = d_0) &= P(\Lambda_{k,j})P(\underline{s}|\underline{y}_s, \delta_k)I[\underline{\theta} \in \Theta_{d_R}] \\ &= L_{\underline{\theta}}(D_0 = d_0). \end{aligned} \tag{1}$$

As the likelihood can be separated into two components, one which does not depend on $\underline{\theta}$ and the other which depends on $\underline{\theta}$ only through the reduced data, by the Neyman Factorization Theorem we can conclude that D_R is a sufficient statistic for $\underline{\theta}$.

To show the minimality of the claim, we will make use of the theorem which gives the following. If \underline{X} is a sample with probability mass/density function $f(\underline{X}|\underline{\theta})$ and $T(\underline{X})$ is a function of \underline{X} such that, for any two sample points \underline{x} and \underline{x}^* , the function $k(f(\underline{x}), f(\underline{x}^*)) > 0$ is constant as a function of $\underline{\theta}$ if and only if $T(\underline{x}) = T(\underline{x}^*)$, then $T(\underline{X})$ is the minimal sufficient statistic for $\underline{\theta}$ (Arnold, 1990; Casella and Berger, 2002).

Take any $d_0 = (\Lambda_{k,j}, (i, y_i) : i \in \underline{s})$ and $d_0^* = (\Lambda_{k^*,j^*}, (i, y_i) : i \in \underline{s}^*)$ where $P(d_0), P(d_0^*) > 0$. Suppose that $P_{\underline{\theta}}(d_0) = k(d_0, d_0^*)P_{\underline{\theta}}(d_0^*)$ where k is independent of $\underline{\theta}$. We can re-express this as

$$P(\Lambda_{k,j})P(\underline{s}|\underline{y}_{\underline{s}}, \delta_k)I[\underline{\theta} \in \Theta_{d_R}] = k(d_0, d_0^*)P(\Lambda_{k^*,j^*})P(\underline{s}^*|\underline{y}_{\underline{s}^*}, \delta_{k^*})I[\underline{\theta} \in \Theta_{d_R^*}] \quad (2)$$

where $d_R = \{(i, y_i) : i \in s\}$, $d_R^* = \{(i, y_i) : i \in s^*\}$. As

$$P(\Lambda_{k,j}), P(\Lambda_{k^*,j^*}), P(\underline{s}|\underline{y}_{\underline{s}}, \delta_k), P(\underline{s}^*|\underline{y}_{\underline{s}^*}, \delta_{k^*}) > 0, \quad (3)$$

the indicators must take on a value of zero or one at the same time. Hence, $I[\underline{\theta} \in \Theta_{d_R}] = I[\underline{\theta} \in \Theta_{d_R^*}]$ for all $\underline{\theta}$ and so $\Theta_{d_R} = \Theta_{d_R^*}$, which implies that $d_R = d_R^*$. Therefore, D_R is the minimal sufficient statistic.

□

Theorem: The random variable D_R is not complete.

Proof: By way of counterexample, consider the case where only one strategy is considered. Then this is the usual case in the classical fixed population approach since mentioning the strategy in the observed data is redundant. Hence, by Cassel et al. (1977) it must be that D_R is not complete.

□

4 Estimation

In this section we derive the expression for the improved point estimators and an estimate of the variance of this estimator. We conclude the section by highlighting how the suggested method of estimation for the variance of the improved estimator avoids the need for estimation of covariances.

4.1 Point estimation

Suppose $\phi = \phi(\underline{\theta})$ is a population quantity to be estimated. Define \mathcal{S} to be the set of all subsets of U that have a positive probability of being obtained with at least one $\delta_k \in \mathcal{D}$. For ease of presentation, we will define

$$\hat{\phi}_{k,j}(s) = \frac{\sum_{\underline{s} \in \mathcal{S}_s} \hat{\phi}_{k,j}(\underline{s}) P(\underline{s} | \underline{y}_{\underline{s}}, \delta_k)}{\sum_{\underline{s} \in \mathcal{S}_s} P(\underline{s} | \underline{y}_{\underline{s}}, \delta_k)} \quad (4)$$

where \mathcal{S}_s is the set of all full samples \underline{s} such that the corresponding reduced set is s . Notice that this is the usual Rao-Blackwellization scheme that is used in traditional survey sampling inference.

Recall that we defined $\Lambda_{k,j}$ to be the strategy that depends on sampling design δ_k and estimator $\hat{\phi}_{k,j}$. With respect to any specific sample s , in the event that s cannot be obtained with δ_k , that is $P(s | \underline{y}_s, \delta_k) = 0$, we will define $\hat{\phi}_{k,j}(s)$ to be zero.

Let $\hat{\phi}(s)$ denote the preliminary estimate of ϕ (preliminary in the sense that improvement has not been made with respect to averaging over estimators corresponding with multiple strategies). The Rao-Blackwellized version of $\hat{\phi}(s)$ is

$$\begin{aligned} \hat{\phi}_{RB}(s) &= E[\hat{\phi}(s) | d_R] \\ &= \frac{\sum_k \sum_j (\hat{\phi}_{k,j}(s) P(\Lambda_{k,j}) P(s | \underline{y}_s, \delta_k))}{\sum_k \sum_j (P(\Lambda_{k,j}) P(s | \underline{y}_s, \delta_k))} \\ &= \frac{\sum_k \sum_j (\hat{\phi}_{k,j}(s) P(\Lambda_{k,j}) P(s | \underline{y}_s, \delta_k))}{P(s | \underline{y}_s)}. \end{aligned} \quad (5)$$

4.2 Variance estimation

The variance of the preliminary estimate $\hat{\phi}$ is

$$\text{var}(\hat{\phi}) = \sum_{s \in \mathcal{S}} \sum_k \sum_j ((\hat{\phi}_{k,j}(s) - E[\hat{\phi}])^2 P(\Lambda_{k,j}) P(s | \underline{y}_s, \delta_k)). \quad (6)$$

The variance of the improved estimator can be determined with the decomposition of variances

as follows,

$$\begin{aligned}
\text{var}(\hat{\phi}_{RB}) &= \text{var}(\hat{\phi}) - E[\text{var}(\hat{\phi}|d_R)] \\
&= \sum_{s \in \mathcal{S}} \sum_k \sum_j ((\hat{\phi}_{k,j}(s) - E[\hat{\phi}])^2 P(\Lambda_{k,j}) P(s|\underline{y}_s, \delta_k)) \\
&\quad - \sum_{s \in \mathcal{S}} \sum_k \sum_j ((\hat{\phi}_{k,j}(s) - \hat{\phi}_{RB}(s))^2 P(\Lambda_{k,j}) P(s|\underline{y}_s, \delta_k)) \\
&= \sum_{s \in \mathcal{S}} \sum_k \sum_j ((\hat{\phi}_{RB}(s) - E[\hat{\phi}])^2 P(\Lambda_{k,j}) P(s|\underline{y}_s, \delta_k)) \\
&= \sum_{s \in \mathcal{S}} ((\hat{\phi}_{RB}(s) - E[\hat{\phi}])^2 P(s|\underline{y}_s)).
\end{aligned} \tag{7}$$

To estimate the variance of the improved estimate, the following estimator can be used. First, an estimate of $\text{var}(\hat{\phi})$ is

$$E[\hat{\text{var}}(\hat{\phi}(s))|d_R] = \sum_k \sum_j (\hat{\text{var}}(\hat{\phi}_{k,j}(s)) P(\Lambda_{k,j}) P(s|\underline{y}_s, \delta_k)) / P(s|\underline{y}_s) \tag{8}$$

where $\hat{\text{var}}(\hat{\phi}_{k,j}(s))$ is the corresponding estimate of the variance of the estimate of $\hat{\phi}_{k,j}(s)$. Notice that this is the Rao-Blackwellized estimate of the preliminary estimate of the variance of $\hat{\phi}_{k,j}(s)$.

Second, an estimate of $E[\text{var}(\hat{\phi}|d_R)]$ is

$$\text{var}(\hat{\phi}(s)|d_R) = \sum_k \sum_j ((\hat{\phi}_{k,j}(s) - \hat{\phi}_{RB}(s))^2 P(\Lambda_{k,j}) P(s|\underline{y}_s, \delta_k)) / P(s|\underline{y}_s). \tag{9}$$

Finally, the estimate of $\hat{\text{var}}(\hat{\phi}_{RB}) = E[\hat{\text{var}}(\hat{\phi}(s))|d_R] - \text{var}(\hat{\phi}(s)|d_R)$ can be used as an estimate of $\text{var}(\hat{\phi}_{RB})$. Furthermore, if $\hat{\text{var}}(\hat{\phi}_{k,j}(s))$ is an unbiased estimate of $\text{var}(\hat{\phi}_{k,j}(s))$ for all k and j , then this estimator is unbiased. We make two remarks here. First, a direct attempt to determining the variance of the improved estimator could be challenging as this would require determining the covariance of estimators corresponding with multiple strategies. In contrast, the aforementioned estimator has the advantage in that it avoids covariance measures of estimators. Second, in some cases this estimator will give negative estimates. In such a case, a conservative approach is to take $\hat{\text{var}}(\hat{\phi}_{RB}) = E[\hat{\text{var}}(\hat{\phi}(s))|d_R]$. However, this is still an outstanding issue and will require future attention.

5 Simulation Studies

In this section we demonstrate how the new inference procedure can assist in various aspects of a survey sampling study. Results from two simulation studies are provided.

5.1 Simulation study 1

The primary objective of this simulation study is to demonstrate how the new inference procedure can allow one to average over estimates corresponding with several strategies to avoid committing to, and possibly obtain a more efficient estimator than that based on, any one of the candidate strategies. We consider the ‘influenza’ data set in R (Manitz et al., 2013) where a count of inhabitants and influenza cases within 424 districts of Germany in 2007 are provided. We aim to estimate the total number of cases of influenza. Two strategies are considered. The first strategy, denoted $\Lambda_{1,1}$, entails the use of a simple random sampling with replacement (SR-SWR) design and the full sample mean. The second strategy, denoted $\Lambda_{2,1}$, entails the use of a probability proportional to size with replacement (PPSWR) sampling design and a Hansen-Hurwitz (HH) estimator where the count of inhabitants serves as the auxiliary information. As we are interested in how the new setup can improve on estimation, Rao-Blackwellization is restricted to directly averaging over the aforementioned estimators corresponding with the two strategies and not with respect to the unique sample elements¹.

$P(\Lambda_{1,1})$ is set to 0.10 and $P(\Lambda_{2,1})$ is set to 0.90 to reflect a level of certainty in the use of a PPSWR design. Samples of size five are obtained and a large number of samples are drawn in order to eliminate any Monte carlo error. Table 1 provides the standardized variance scores and corresponding coverage rates based on nominal 95% confidence intervals via the Central Limit Theorem (CLT).

¹Though sampling is carried out with replacement, we consider small sample sizes relative to the population size and hence one may argue that the use of the Hansen-Hurwitz estimator is practical for such a case. Further, one can avoid the tedious computations required for inclusion probabilities typically required for without-replacement estimators.

Table 1: Simulation study based on the influenza data set in R. The sample size is five and $P(\Lambda_{1,1}) = 0.10, P(\Lambda_{2,1}) = 0.90$. Population total is 18,900. All estimators are unbiased. The variance scores are standardized about the Hansen-Hurwitz estimator used with the PPSWR design (that is, those obtained with strategy $\Lambda_{2,1}$).

Strategy	Standardized variance scores	Coverage rates
$\Lambda_{1,1}$	1.330	0.971
$\Lambda_{2,1}$	1	0.886
Preliminary	1.033	0.895
Improved	0.961	0.913

In this case the improved estimator outperforms each of the individual strategy-based and preliminary estimators. Further, the coverage rates of the improved estimator rest between those based on the two strategies. To further investigate the details of the study, Figure 1 provides a plot of the preliminary by improved estimators for 2000 randomly selected points. The sampling design used to select the sample is reflected upon by the character of the plot. Notice how the improved estimators are typically pulled towards their expectation, and in general there is a dampening effect on the preliminary estimates with Rao-Blackwellization. This is especially evident when such estimators are extreme. In some cases averaging over the two strategies is intuitive since the sample may *appear* as if it were selected under a design corresponding with a strategy not chosen at the initial stage of the study (in particular when extreme preliminary estimates are reported). Note that if considering designs like SRSWR and PPSWR when sample sizes are small, this could very easily be the case.

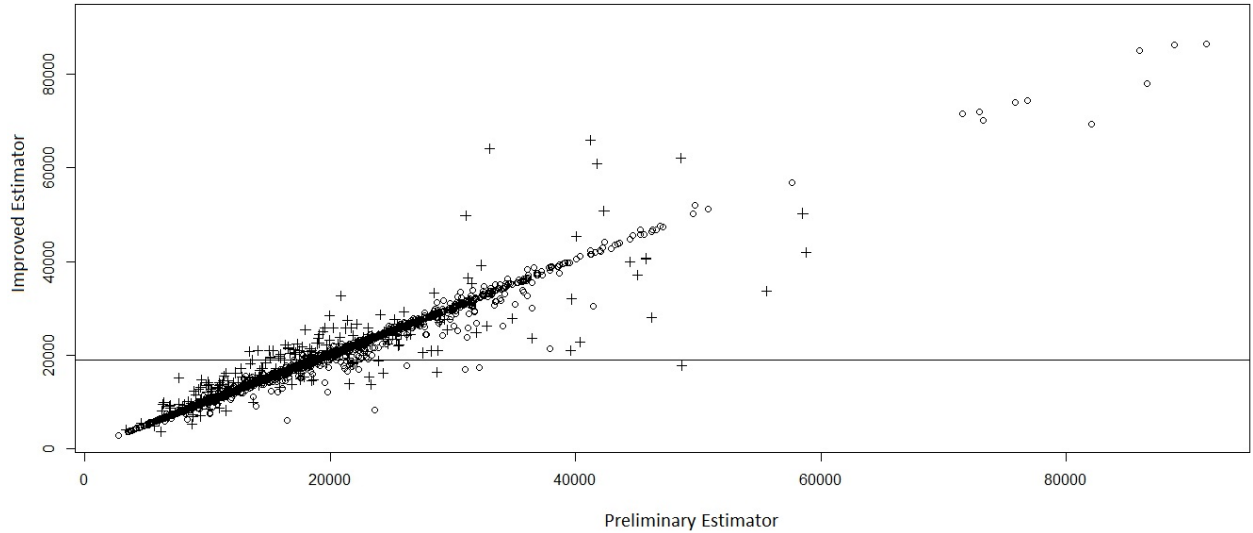


Figure 1: A plot of 2000 randomly selected points from the simulation study. Preliminary estimates correspond with the x-axis and improved estimates correspond with the y-axis. The crosses correspond with the use of $\Lambda_{1,1}$ and circles correspond with the use of $\Lambda_{2,1}$. The vertical and horizontal lines correspond with the true population total.

The simulation study is extended over a range of sample sizes. It is found that with sizes less than or equal to fifteen, the performance of the improved estimator is noticeably better than each of the strategy-based and preliminary estimators. Hence, we can conclude that in some small sample studies the use of the inference procedure we have outlined may serve as a competitive alternative to, and may result in a more efficient than some of, commonly used estimators based on the traditional setup. For sample sizes greater than fifteen the improved estimator is weighted almost entirely towards the strategy used at the preliminary stage. Hence the benefit from using the new inference procedure in these cases is negligible. This is primarily due to the probability of selecting the sample under the design corresponding with the strategy chosen at the initial stage of the study relative to the probability of selecting it with the design corresponding with the strategy that is not chosen.

5.2 Simulation study 2

The primary objective of this simulation study is to demonstrate how the new inference procedure can assist in a retrospective analysis. The details of the study are provided below.

We consider a case where the sample is selected via SRSWR but the sample size is unknown; only the unique elements of the sample are recorded. Features of the effective sample mean are

of interest, in particular a variance estimator of the sample mean. We apply our inference setup to obtain an improved estimate of the variance of the effective sample mean estimator.

The effective sample mean is defined to be $\bar{y}_\nu = \frac{1}{\nu} \sum_{i \in s} y_i$ where ν is the size of the unique sample s . The variance of the effective sample mean estimator (Pathak, 1962) is

$$\text{var}(\bar{y}_\nu) = \frac{\sum_{j=1}^{N-1} j^{n-1}}{N^n} \frac{1}{N-1} \sum_{i=1}^N (Y_i - \mu)^2 \quad (10)$$

where μ is the population mean and n is the number of draws made to select the sample. An unbiased estimator for $\text{var}(\bar{y}_\nu)$ is

$$\hat{\text{var}}(\bar{y}_\nu) = \frac{\sum_{j=1}^{N-1} j^{n-1}}{N^n} \frac{1}{\nu-1} \sum_{i \in s} (Y_i - \bar{y}_\nu)^2. \quad (11)$$

We shall define δ_k to be the SRSWR design that depends on a specific final sample size where repeat selections may occur, n_k say. We shall define $\hat{\phi}_{k,1}$ to be the estimator in (11) that depends on the corresponding sample size associated with δ_k . Hence, the improved version of the estimator presented in (11) is of the form

$$\begin{aligned} \hat{\text{var}}(\bar{y}_\nu)_{RB} &= E[\hat{\text{var}}(\bar{y}_\nu)_{n_k}(s) | d_R] \\ &= \frac{\sum_k (\hat{\text{var}}(\bar{y}_\nu)_{n_k}(s) P(\Lambda_{n_k}) P(s | \delta_k))}{P(s)}. \end{aligned} \quad (12)$$

where $P(\Lambda_{n_k})$ is the probability that the sample size is n_k and $\hat{\text{var}}(\bar{y}_\nu)_{n_k}(s)$ is the estimate of the variance of \bar{y}_ν when the sample size is n_k . Notice that this estimator will depend on evaluating the probability of selecting the unique sample under a SRSWR design for a series of specific final sample sizes. We have not come across an expression in the literature for evaluating such a probability and have therefore taken the liberty of providing its derivation in the Appendix.

In our study we consider the ‘trees’ dataset (R Core Team, 2013), a data set that provides measurements of the girth, height, and volume of timber in 31 felled black cherry trees. We define $\Lambda_{k,1}$ to be based on a SRSWR design with a final sample size of $5 + k - 1$ where $k = 1, 2, \dots, 10$. In each case, $\hat{\phi}_{k,1}$ is the estimate of the variance of the effective sample mean based on the corresponding sample size of $5 + k - 1$. In the simulation study strategies are

chosen such that $P(\Lambda_{k,1}) \propto \frac{1}{k}$ so that smaller sample sizes are expected. Table 2 summarizes the results based on a large number of simulation runs. The naive estimator is based on the actual effective sample size (that is, n is treated as ν so that inference is based on a hypothetical sample that is equivalent to the unique sample). Coverage results are based on nominal 95% CLT confidence intervals.

Table 2: Simulation study 2 results. Results based on the girth, height, and volume responses, respectively.

Response	Estimator	Relative Bias	Variance	Coverage
Girth				
	Naive	0.138	0.783	0.904
	Preliminary	0	0.593	0.889
	Improved	0	0.563	0.890
Volume				
	Naive	0.131	13.00	0.911
	Preliminary	0	9.94	0.894
	Improved	0	9.46	0.895
Height				
	Naive	0.131	826	0.882
	Preliminary	0	626	0.866
	Improved	0	601	0.868

In each case, the new inference procedure provides improved estimates, and therefore tighter confidence bands corresponding with the estimates for the population mean, as well as coverage rates of the population mean that are on par with those found using the preliminary estimator. Hence, we can conclude that the new inferential setup has the potential to lend itself well for obtaining efficient estimators for sampling distributions when only limited information pertaining to the sampling design is available.

6 Implications

In this manuscript we have introduced a new approach to making design-based inference in survey studies; the framework we outline rests on incorporating an additional element of randomization via allowing for the strategy to first be selected at random from amongst a set of candidate strategies. For such a case, we have derived the minimal sufficient statistic for the

population parameter vector and have shown that it is not complete. We demonstrate how the new approach can be useful for some specific survey sampling studies.

In general, the new inference framework allows the analyst to make a compromise over a set of candidate strategies through the Rao-Blackwellized estimator and is therefore not required to fully commit to one potentially weak strategy. In the first simulation study we demonstrate that studies based on small sample sizes can benefit from basing a final estimator on an average of estimators corresponding with strategies that are comprised of various sampling designs. As we highlighted, this approach works especially well when extreme values are reported for preliminary estimators. Research on how small sample studies that usually encounter such occurrences can benefit from this approach will make for useful future work.

In an empirical setting the analyst may not have full control over the sample selection procedure. Therefore, they may only have a limited amount of information pertaining to the sampling design that was used for data collection. An example motivates the second simulation study; when sampling is carried out via a with replacement sampling design it could be the case that only the unique elements of the sample are provided to the analyst, perhaps due to a lack of communication. In such occurrences one can posit a series of sample sizes and determine Rao-Blackwell-type estimators for population quantities. We further explored this approach for the simulation parameters outlined in the second simulation study. We consistently sampled at a fixed final sample size and posited sample sizes that ranged up to twice the final sample size. We also posited equal probabilities for the sample sizes. We found that coverage rates corresponding with the improved estimator were only negligibly affected while the variance was significantly reduced.

The approach presented in this paper suggests a new method for applying a Bayes like analysis to a sampling setup; consider Expression (5) and how it is reflective of a Bayes estimator. In the case of when limited knowledge is available on the sampling design used, one can posit a series of (prior) sampling designs with a corresponding (prior) distribution that each design was used and then base inference on the Rao-Blackwell expression. Future work on combinations of strategies that can be used to draw meaningful conclusions in certain situations when the analyst has limited control/knowledge of the sample selection procedure is deserving of future attention.

References

- Arnold, S. (1990). *Mathematical Statistics*. Englewood Cliffs, New Jersey: Prentice-Hall.
- Basu, D. (1969). Role of the sufficiency and likelihood principles in sample survey theory. *The Indian Journal of Statistics, Series A* **31**, 441–454.
- Casella, G. and Berger, R. L. (2002). *Statistical Inference*. Wadsworth and Brooks/Cole, second edition.
- Cassel, C.-M., Särndal, C.-E., and Wretman, J. H. (1977). *Foundations of inference in survey sampling*. Wiley New York.
- Godambe, V. (1966). A new approach to sampling from finite populations i: Sufficiency and linear estimation. *Journal of the Royal Statistical Society. Series B (Methodological)* **28**, 310–319.
- Manitz, J., contributions by Mark Hempelmann, Kauermann, G., Kuechenhoff, H., Shao, S., Oberhauser, C., Westerheide, N., and Wiesenfarth, M. (2013). *samplingbook: Survey Sampling Procedures*. R package version 1.2.0.
- Pathak, P. K. (1962). On simple random sampling with replacement. *Sankhya: The Indian Journal of Statistics, Series A (1961-2002)* **24**, pp. 287–302.
- Pathak, P. K. (1964). Sufficiency in sampling theory. *The Annals of Mathematical Statistics* **35**, pp. 795–808.
- R Core Team (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Thompson, S. K. and Seber, G. A. F. (1996). *Adaptive Sampling*. Wiley Series in Probability and Statistics, New York.

A Appendix

In this appendix we derive the probability of obtaining a specific sample of unique elements where the size of the unique sample is ν , the size of the full sample is n , and the population size is N .

Suppose that units are selected from a population where $U = \{1, 2, \dots, N\}$ are the population unit labels and where selections are made via the simple random sampling with replacement (SRSWR) design of size n . Note that this is equivalent to drawing one unit at a time (where each unit has probability $1/N$ of being selected on each draw), recording the label identification, and placing the unit back in the population. Suppose that \underline{s} is the original ordered sample. Let s be the effective (unique) sample. For example, if $\underline{s} = (2, 3, 1, 1, 4, 2)$ then $s = \{1, 2, 3, 4\}$. We would like to determine the probability of obtaining s , $P(s)$.

Partition U into s and $U \setminus s$. Now, define X to be the event that all n draws are made from s , and define Y to be the event that no units in s are missing from the final sample. We would like to determine $P(s) = P(X \cap Y) = P(X) \times P(Y|X)$. Suppose $|s| = \nu$. Clearly, $P(X) = (\frac{\nu}{N})^n$. Now, define A_i to be the event that unit i is selected for the final sample given that all draws are made from s , $i \in s$. We will now determine

$$P(Y|X) = 1 - P(\cup_{i \in s} \bar{A}_i). \quad (13)$$

By the principle of inclusion-exclusion,

$$P(\cup_{i \in s} \bar{A}_i) = \sum_{k=1}^{\nu} \left((-1)^{k-1} \sum_{I \subset s: |I|=k} P(\bar{A}_I) \right). \quad (14)$$

If $|I| = k$, then $P(\bar{A}_I) = \left(\frac{\nu-k}{\nu} \right)^n$ so that expression (14) can be rewritten as

$$\sum_{k=1}^{\nu} \left((-1)^{k-1} \binom{\nu}{k} \left(\frac{\nu-k}{\nu} \right)^n \right) \quad (15)$$

since there are $\binom{\nu}{k}$ subsets of s that are of size k .

Therefore,

$$P(s) = \left(\frac{\nu}{N} \right)^n \left[1 - \sum_{k=1}^{\nu} \left((-1)^{k-1} \binom{\nu}{k} \left(\frac{\nu-k}{\nu} \right)^n \right) \right]. \quad (16)$$